

AMENDMENTS TO THE CLAIMS

Claims 1-3, 5-10, 14, 15, 17, 19-22, 24, 28, 30-32 and 35 are amended herein.

Claim 11-13, 25-27 are cancelled.

Claims 1-10, 14-24, and 28-37 are now pending. All pending claims are produced below. In addition, the status of each claim is also indicated below and appropriately noted as “Original”, “Currently Amended”, “Canceled”, “New”, “Withdrawn”, “Previously Presented”, and “Not Entered” as requested by the Office.

1. (Currently Amended) A system for identifying language attributes through probabilistic analysis, comprising:

a storage storage system adapted to store a set of language classes, which each identify a language and a character set encoding, and further adapted to store a plurality of training documents;

an attribute modeler adapted to train an attribute model by evaluating occurrences of one or more document properties within each the training document documents and, for each language class, calculating a probability for the document properties set conditioned on the occurrence of the language class, the trained attribute model stored in the storage; and

a text modeler adapted to train a text model by evaluating byte occurrences within each the training document documents and, for each language class, calculating a probability for the byte occurrences conditioned on the occurrence of the language class, the trained text model stored in the storage.

2. (Currently Amended) A system according to Claim 1, further comprising:

a training engine calculating adapted to calculate an overall probability for each language class by evaluating the probability for the document properties set and the probability for the byte occurrences.

3. (Currently Amended) A system according to Claim 1, further comprising:
an assignment module assigning adapted to assign the overall probability for each a
language class in accordance with the formula:

$$\arg \max_{cls} P(text | cls) \cdot P(props | cls) \cdot P(cls)$$

where *cls* is the language class, *text* is the byte occurrences set, *props* are the document properties, and $P(text | cls)$ is the probability for the byte occurrences, and $P(props | cls)$ is the probability for the document properties set.

4. (Original) A system according to Claim 1, wherein the document properties comprise at least one of top level domain, HTTP content character set encoding and language header parameters, and HTML content character set encoding and language metatags.

5. (Currently Amended) A system according to Claim 4, further comprising:
an assignment module assigning adapted to assign the probability for the document properties set based on the attribute model in accordance with the formula:

$$P(tld, enc | cls) \cdot P(cls)$$

where *tld* is the top level domain, *enc* is the character set encoding and *cls* is the language class.

6. (Currently Amended) A system according to Claim 1, further comprising:
a counting module counting adapted to count byte co-occurrences within each a
training document, and determining determine the probability for the byte occurrences based on the byte co-occurrences.

7. (Currently Amended) A system according to Claim 6, wherein the byte co-occurrences comprise a set of trigrams, further comprising:

a probability module calculating adapted to calculate a probability of each a trigram as the number of occurrences of the trigram divided by the total number of trigram occurrences in each of the training documents for each a language class.

8. (Currently Amended) A system according to Claim 7, further comprising:
an assignment module assigning adapted to assign the probability for the byte
occurrences set based on the text model in accordance with the formula:

$$P(text | cls)$$

where *text* is the set of trigrams and *cls* is the language class.

9. (Currently Amended) A system according to Claim 1, further comprising:
a training engine performing adapted to perform iterative training by providing the
probability for the document properties set and the probability for the byte
occurrences set respectively to the evaluation of byte occurrences and
assignment of the set of language classes.

10. (Currently Amended) A system according to Claim 1, further comprising:
a back off module evaluating adapted to evaluate less frequently occurring document
properties by calculating a probability for each a less frequently occurring
document property conditioned on the occurrence of the language class.

11. (Canceled)

12. (Canceled)

13. (Canceled)

14. (Currently Amended) A system according to Claim 1, wherein each at least one
training document comprises one of a Web page and a news message.

15. (Currently Amended) A method for identifying language attributes through
probabilistic analysis, comprising:

defining a set of language classes, which each identify a language and a character set
encoding, and a plurality of training documents;
evaluating occurrences of one or more document properties within each the training
document documents and, for each language class, calculating a probability

for the document properties set conditioned on the occurrence of the language class by an attribute model; and

evaluating byte occurrences within ~~each~~ the training ~~document documents~~ and, for each language class, calculating a probability for the byte occurrences conditioned on the occurrence of the language class by a text model.

16. (Original) A method according to Claim 15, further comprising:

calculating an overall probability for each language class by evaluating the probability for the document properties set and the probability for the byte occurrences.

17. (Currently Amended) A method according to Claim 15, further comprising:

assigning the overall probability for ~~each~~ a language class in accordance with the formula:

$$\arg \max_{cls} P(text | cls) \cdot P(props | cls) \cdot P(cls)$$

where *cls* is the language class, *text* is the byte occurrences set, *props* are the document properties, and $P(text | cls)$ is the probability for the byte occurrences, and $P(props | cls)$ is the probability for the document properties set.

18. (Original) A method according to Claim 15, wherein the document properties comprise at least one of top level domain, HTTP content character set encoding and language header parameters, and HTML content character set encoding and language metatags.

19. (Currently Amended) A method according to Claim 18, further comprising:

assigning the probability for the document properties set based on the attribute model in accordance with the formula:

$$P(tld, enc | cls) \cdot P(cls)$$

where *tld* is the top level domain, *enc* is the character set encoding and *cls* is the language class.

20. (Currently Amended) A method according to Claim 15, further comprising:
counting byte co-occurrences within ~~each~~ a training document; and
determining the probability for the byte occurrences based on the byte co-
occurrences.

21. (Currently Amended) A method according to Claim 20, wherein the byte co-
occurrences comprise a set of trigrams, further comprising:
calculating a probability of ~~each~~ a trigram as the number of occurrences of the trigram
divided by the total number of trigram occurrences in ~~each~~ of the training
documents for ~~each~~ a language class.

22. (Currently Amended) A method according to Claim 21, further comprising:
assigning the probability for the byte occurrences set based on the text model in
accordance with the formula:

$$P(text | cls)$$

where *text* is the set of trigrams and *cls* is the language class.

23. (Original) A method according to Claim 15, further comprising:
performing iterative training by providing the probability for the document properties
set and the probability for the byte occurrences set respectively to the
evaluation of byte occurrences and assignment of the set of language classes.

24. (Currently Amended) A method according to Claim 15, further comprising:
evaluating less frequently occurring document properties by calculating a probability
for ~~each~~ a less frequently occurring document property conditioned on the
occurrence of the language class.

25. (Canceled)
26. (Canceled).
27. (Canceled)

28. (Currently Amended) A method according to Claim 15, wherein ~~each~~ at least one training document comprises one of a Web page and a news message.

29. (Original) A computer-readable storage medium holding code for performing the method according to Claim 15.

30. (Currently Amended) A system for identifying documents by language using probabilistic analysis of language attributes, comprising:

a set of language classes, each language class comprising a language name and a character set encoding name;

a training corpora comprising a plurality of training documents;

an attribute modeler ~~training adapted to train~~ an attribute model by evaluating a top level domain and character set encoding associated with ~~each~~ the training ~~document documents~~ and, for each language class, calculating a probability for each such top level domain and character set encoding conditioned on the occurrence of the each language class; and

a text modeler ~~training adapted to train~~ a text model by evaluating co-occurrences of a plurality of bytes within ~~each~~ the training ~~document documents~~ and, for each language class, calculating a probability for the byte co-occurrences conditioned on the occurrence of the each language class.

31. (Currently Amended) A system according to Claim 30, further comprising:

a training engine ~~calculating adapted to calculate~~ an overall probability for each language class by evaluating the probability for the top level domain and character set encoding based on the attribute model and the probability for the byte occurrences based on the text model.

32. (Currently Amended) A system according to Claim 31, further comprising:

~~a classifier classifying one or more documents, comprising:~~

~~an attribute evaluator evaluating a top level domain and character set encoding in each document and applying the attribute model to the evaluated top level domain and character set encoding;~~

~~a text evaluator evaluating byte occurrences in each document and applying the text model to the evaluated byte occurrences; and~~

~~an assignment module assigning at least one language class based on the applications of the attribute model and the text model.~~

a plurality of unlabeled documents; and

a classifier classifying one or more unlabeled documents by at least one language class, comprising:

an attribute evaluator determining document properties within the documents and initializing language class probability to each document from the attribute model;

a text evaluator evaluating byte occurrences in the documents and updating the language class probability of the each document from the text model;

a pruner pruning at least one language class falling below a predetermined probability threshold; and

an assignment module assigning at least one language class based on the language class probability of each document.

33. (Original) A method for identifying documents by language using probabilistic analysis of language attributes, comprising:

defining a set of language classes, each language class comprising a language name and a character set encoding name;

assembling a training corpora comprising a plurality of training documents;

training an attribute model by evaluating a top level domain and character set encoding associated with each training document and, for each language class, calculating a probability for each such top level domain and character set encoding conditioned on the occurrence of the each language class; and

training a text model by evaluating co-occurrences of a plurality of bytes within each training document and, for each language class, calculating a probability for

the byte co-occurrences conditioned on the occurrence of the each language class.

34. (Original) A method according to Claim 33, further comprising:

calculating an overall probability for each language class by evaluating the probability for the top level domain and character set encoding based on the attribute model and the probability for the byte occurrences based on the text model.

35. (Currently Amended) A method according to Claim [[34]] 33, further comprising:

~~classifying one or more documents, comprising:~~

~~evaluating a top level domain and character set encoding in each document and applying the attribute model to the evaluated top level domain and character set encoding;~~

~~evaluating byte occurrences in each document and applying the text model to the evaluated byte occurrences; and~~

~~assigning at least one language class based on the applications of the attribute model and the text model~~

accessing a plurality of unlabeled documents; and

classifying one or more unlabeled documents by at least one language class, comprising:

determining document properties within the documents and initializing language class probability to each document from the attribute model;

evaluating byte occurrences in each document and updating the language class probability of the document from the text model;

pruning at least one language class failing below a predetermined probability threshold; and

assigning at least one language class based on the language class probability of the document.

36. (Original) A computer-readable storage medium holding code for performing the method according to Claim 30.

37. (Original) An apparatus for identifying documents by language using probabilistic analysis of language attributes, comprising:

means for defining a set of language classes, each language class comprising a language name and a character set encoding name;

means for training an attribute model by assigning at least one top level domain and character set encoding pairing to at least one language class for each of a plurality of training documents and calculating a probability for each such top level domain and character set encoding pairing conditioned on the occurrence of the assigned language class; and

means for training a text model by evaluating co-occurrences of a plurality of bytes within each training document and, for each language class, calculating a probability for the byte co-occurrences conditioned on the occurrence of the language class based on the attribute model.